

**APPLICATION FOR
UNITED STATES PATENT
IN THE NAME**

Of

KENJI YAMAGAMI

For

**A FAILURE NOTIFICATION METHOD AND SYSTEM USING
REMOTE MIRRORING FOR CLUSTERING SYSTEMS**

DOCKET NO. 36992.00068

**Please direct communications to:
SQUIRE, SANDERS & DEMPSEY L.L.P.
600 Hansen Way
Palo Alto, CA 94304-1043
(650) 856-6500
Express Mail Number: EL701362040US**

09760345.01201
T02T0"54E09/60

1 A FAILURE NOTIFICATION METHOD AND SYSTEM
2 USING REMOTE MIRRORING FOR CLUSTERING SYSTEMS
3

4 Inventor: KENJI YAMAGAMI
5

6 FIELD OF THE INVENTION

7 The present invention relates to cluster computing systems,
8 and relates more particularly to systems and methods for
9 providing heartbeat-checking mechanisms by use of remote mirror
10 technology for cluster computing systems. The present invention
11 permits a host on a primary site to send heartbeat signals to a
12 host on a secondary site (and vice versa) by use of remote mirror
13 technology.
14

15 BACKGROUND OF THE INVENTION

16 "Clustering" is the known technique of connecting multiple
17 computers (or host servers) and enabling the connected computers
18 to act like a single machine. Clustering is used for parallel
19 processing, for load balancing, and for fault tolerance.
20 Corporations often cluster servers together in order to
21 distribute computing-intensive tasks and risks. If one server in
22 the cluster computing system fails, then an operating system can
23 move its processes to a non-failing server in the cluster
24 computing system, and this allows end users to continue working
25 while the failing server is revived.

Cluster computing systems are becoming popular for preventing operation interruptions of applications. Some cluster computing systems have two groups of hosts (e.g., servers), wherein one host group works as the production system, while the other host group works as the standby system. One host group is typically geographically dispersed (e.g., several hundred miles) from the other host group. Each host group has its own associated storage system (e.g., a disk system). These two storage systems typically implement remote mirroring technology which is discussed below. Therefore, the associated storage system connecting to the standby host group contains the same data as the associated storage system connecting to the production host group.

The network connecting two host server groups is typically a Wide Area Network (WAN), such as the Internet. WANs are not typically reliable since WANs are often subject to failure. Transfer of data across the Internet can be subject to delays and can lead to data loss. Therefore, the standby host group may inappropriately take over the processes of the production host group (even if there is no failure in the production host group) because the standby host group may erroneously see a network problem (e.g., link failure or data transmission delay) as a failure state of the production host group.

1 The host group in the production system may access a storage
2 volume commonly known a primary volume (PVOL) in the associated
3 storage system of the production system host group. Similarly,
4 the host group in the standby system may access a storage volume
5 commonly known a secondary volume (SVOL) in the associated
6 storage system of the standby system host group. The primary
7 volume (PVOL) is mirrored by the secondary volume (SVOL). A
8 storage system may have both PVOLs and SVOLs.

9 Storage-based remote mirroring technology creates and stores
10 mirrored volumes of data between a given distance. Two disk
11 systems are directly connected by remote links such as an
12 Enterprise System Connectivity (ESCON) architecture, Fibre
13 Channel, telecommunication lines, or a combination of these
14 remote links. The data in the local disk system is transmitted
15 to (via the remote links) and copied in the remote disk system.
16 These remote links are typically highly reliable, in comparison
17 to a usual network such as the Internet. If an unreliable remote
18 link fails, then this failure may disadvantageously result in the
19 loss of data.

20 U.S. Patent Nos. 5,459,857 and 5,544,347 both disclose
21 remote mirroring technology. These patent references disclose
22 two disk systems connected by remote links, with the two disk
23 systems separated by a distance. Mirrored data is stored in
24 disks in the local disk system and in the remote disk system.

05760345-011201

1 The local disk system copies data on a local disk when pair
2 creation is indicated. When a host server updates data on the
3 disk, the local disk system transfers the data to the remote disk
4 system through the remote link. Thus, host operation is not
5 required to maintain a mirror data image of one disk system in
6 another disk system.

7 U.S. Patent No. 5,933,653 discloses another type of data
8 transferring method between a local disk system and a remote disk
9 system. In the synchronous mode, the local disk system transfers
10 data to the remote disk system before completing a write request
11 from a host. In the semi-synchronous mode, the local disk system
12 completes a write request from the host and then transfers the
13 write data to the remote disk system. Subsequent write requests
14 from the host are not processed until the local disk system
15 completes the transfer of the previous data to the remote disk
16 system. In the adaptive copy mode, pending data to be
17 transferred to the remote disk system is stored in a memory and
18 transferred to the remote disk system when the local disk system
19 and/or remote links are available for the copy task.

20 There is a need for a system and method that will overcome
21 the above-mentioned deficiencies of conventional methods and
22 systems. There is also a need for a system and method that will
23 increase reliability of cluster computing systems and improved
24 failure detection in these computing systems. There is also a

1 need for a system and method that will accurately detect failure
2 in the production host group of a cluster system so that the
3 standby host group is prevented from taking over the processes of
4 the production host group when the production host group has not
5 failed.

6

7

05760345 011201

1 SUMMARY

2 The apparatus and methods described in this invention
3 provide heartbeat-checking mechanisms by using remote mirror
4 technology for cluster computing systems. Once the remote
5 mirrors are created and set up for heartbeat checking functions,
6 a first host sends heartbeat message to another host that is
7 geographically dispersed from the first host. The heartbeat
8 signals are transmitted through a network and/or by use of
9 remote mirrors.

10 In one embodiment, the present invention broadly provides a
11 cluster computing system, comprising: a production host group; a
12 standby host group coupled to the production host group by a
13 network; and a remote mirror coupled between the production host
14 group and the standby host group, the remote mirror including a
15 production site heartbeat storage volume (heartbeat PVOL) and a
16 standby site heartbeat storage volume (heartbeat SVOL) coupled
17 by a remote link to the heartbeat PVOL, with the production host
18 group configured to selectively send a heartbeat signal to the
19 standby host group by use of at least one of the network and the
20 remote link.

21 In another embodiment, the present invention enables the
22 bi-directional transmission of a heartbeat signal. The cluster
23 computing system may comprise a second remote mirror coupled
24 between the production host group and the standby host group,

1 the second remote mirror including a second remote link for
2 transmitting a heartbeat signal, and the standby host group is
3 configured to selectively send a heartbeat signal to the
4 production host group by use of at least one of the network and
5 the second remote link.

6 In another embodiment, the present invention broadly
7 provides a method of checking for failure in a cluster computing
8 system. The method comprises: generating a heartbeat signal
9 from a production host group; selectively sending the heartbeat
10 signal to the standby host group from the production host group
11 by use of at least one of a network and a remote link; and
12 enabling the standby host group to manage operations of the
13 cluster computing system if an invalid heartbeat signal is
14 received by the standby host group from the production host
15 group.

16 In another embodiment, the present invention provides a
17 method of installing remote mirrors in a cluster computing
18 system. The method comprises: registering a first storage
19 volume to a device address entry, the first storage volume
20 located in a production site; from the production site, changing
21 a remote mirror that includes the first storage volume into an
22 enabled mode; sending an activation message from the production
23 site to a standby site; registering a second storage volume to
24 the device address entry, the second storage volume located in

1 the standby site; and from the standby site, changing the remote
2 mirror into an enabled mode to install a remote mirror formed by
3 the first storage volume and second storage volume.

4 In another embodiment, the present invention provides a
5 method of de-installing remote mirrors in a cluster computing
6 system. The method comprises: from a production site, changing
7 a remote mirror into a disabled mode; sending a de-activation
8 message from the first production site to a standby site; and
9 from the standby site, changing the remote mirror into a
0 disabled mode to de-install the remote mirror.

1 In another embodiment, the present invention provides a
2 method of transmitting a heartbeat message from a production site
3 host to a standby site host in a cluster computing system. The
4 method comprises: determining if a network between the production
5 site host and the standby site host is enabled; if the network is
6 enabled, sending a heartbeat message along the network from the
7 production site host to the standby site host; determining if a
8 remote mirror between the production site host and the standby
9 site host is enabled; and if the remote mirror is enabled,
0 sending a heartbeat message along the remote mirror from the
1 production site host to the standby site host.

2 In another embodiment, the present invention provides a
3 method of receiving a heartbeat message from a production site
4 host to a standby site host in a cluster computing system. The

1 method comprises: determining if a network between the production
2 site host and the standby site host is enabled; if the network is
3 enabled, checking for a heartbeat message along the network from
4 the production site host to the standby site host; determining if
5 a remote mirror between the production site host and the standby
6 site host is enabled; if the remote mirror is enabled, checking
7 for a heartbeat message along the remote mirror from the
8 production site host to the standby site host; and if an invalid
9 heartbeat is received along the network and along the remote
10 mirror, enabling the standby host to manage operations of the
11 cluster computing system.

12 In another embodiment, the present invention provides a
13 method of setting a heartbeat checking procedure between a
14 primary group and a secondary group in a cluster computing
15 system. The method comprises: providing a request command that
16 determines the heartbeat checking procedure; responsive to the
17 request command, enabling a first heartbeat check module in the
18 primary group to activate or de-activate a network between the
19 primary group and the secondary group; responsive to the request
20 command, enabling the first heartbeat check module to activate
21 or de-activate a remote mirror between the primary group and the
22 secondary group; permitting the first heartbeat check module to
23 send the request command to a second heartbeat check module in
24 the secondary group; responsive to the request command, enabling

1 the second heartbeat check module to activate or de-activate the
2 network between the primary group and the secondary group;
3 responsive to the request command, enabling the second heartbeat
4 check module to activate or de-activate the remote mirror
5 between the primary group and the secondary group; if the second
6 heartbeat check module has activated the network, then checking
7 for a heartbeat signal along the network; and if the second
8 heartbeat check module has activated the remote mirror, then
9 checking for a heartbeat signal along the remote mirror.

0 The present invention may advantageously provide a system
1 and method that increase the reliability of cluster computing
2 systems and improve failure detection in these computing systems.
3 The present invention may also advantageously provide a system
4 and method that will accurately detect failure in the production
5 host group of a cluster system so that the standby host group is
6 prevented from taking over the processes of the production host
7 group when the production host group has not failed. The present
8 invention may also advantageously provide a system and method
9 that permit a production host group to check for heartbeat
10 signals from a standby host group.

11

12

13

1 BRIEF DESCRIPTION OF THE DRAWINGS

2 Figure 1 is block diagram of a system configuration in
3 accordance with an embodiment of the present invention;

4 Figure 2 is a block diagram showing an example of a
5 Heartbeat Status Table stored in each of the master hosts shown
6 in Figure 1, in accordance with an embodiment of the present
7 invention;

8 Figure 3 is a block diagram showing an example of the data
9 format of a heartbeat message, in accordance with an embodiment
10 of the present invention;

11 Figure 4 is a flowchart diagram illustrating a method of
12 installing a mirror used for heartbeat signals, in accordance
13 with an embodiment of the present invention;

14 Figure 5 is a flowchart diagram illustrating a method of de-
15 installing a mirror used for heartbeat signals, in accordance
16 with an embodiment of the present invention;

17 Figure 6 is a flowchart diagram illustrating a method of
18 sending a heartbeat message, in accordance with an embodiment of
19 the present invention;

20 Figure 7 is a flowchart diagram illustrating a method of
21 receiving a heartbeat message, in accordance with an embodiment
22 of the present invention;

1 Figure 8 is a flowchart diagram illustrating a method of
2 setting the heartbeat checking procedure, in accordance with an
3 embodiment of the present invention;

4 Figure 9 is block diagram of a system configuration in
5 accordance with another embodiment of the present invention;

6 Figure 10 is a flowchart diagram illustrating a method of
7 failure notification in accordance with an embodiment of the
8 present invention; and

9 Figure 11 is a block diagram illustrating an example of a
0 format of a failure indication message in accordance with an
1 embodiment of the present invention.

009760345 "011201
T02T050345

1 DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS

2 The following description is provided to enable any person
3 skilled in the art to make and use the present invention, and is
4 provided in the context of a particular application and its
5 requirements. Various modifications to the embodiments will be
6 readily apparent to those skilled in the art, and the generic
7 principles defined herein may be applied to other embodiments and
8 applications without departing from the spirit and scope of the
9 present invention. Thus, the present invention is not intended
10 to be limited to the embodiments shown, but is to be accorded the
11 widest scope consistent with the principles, features, and
12 teachings disclosed herein.

13 Figure 1 is a block diagram of a system 50 in accordance
14 with an embodiment of the present invention. The system 50
15 comprises two host groups which are shown as primary group
16 (production host group) 130a and secondary group (standby host
17 group) 130b. The primary group 130a is typically located in a
18 production site and is remote from the secondary group 130b which
19 is typically located in a standby site. The primary group 130a
20 comprises one or more hosts 100a, and the secondary group 130b
21 comprises one or more hosts 100b. The hosts are typically
22 servers.

23 As known to those skilled in the art, a server is a computer
24 or device on a network that manages network resources. For

09760345 "011201
1001

1 example, a file server is a computer and storage device dedicated
2 to storing files. Any user on the network can store files on the
3 server. A print server is a computer that manages one or more
4 printers, and a network server is a computer that manages network
5 traffic. A database server is a computer system that processes
6 database queries.

7 An application 103a normally runs at the primary group 130a,
8 while an application 103b at the secondary group 130b is in the
9 standby mode, as conventionally known in cluster computing
10 systems. When a heartbeat check 101b (in secondary group 130b)
11 determines that the heartbeat check 101a has failed, then
12 application 103a "fails over" to the secondary group 130b in the
13 standby site. In other words, when the application 103a fails
14 over to the secondary group 130b, then the application 103b in
15 the secondary group 130b will run for the system 50.

16 The application 103a will also fail over to the secondary
17 group 130a and the application 103b will run for system 50 when
18 the heartbeat check 101a determines that it is unable to function
19 any longer; this would occur, for example, when only one host
20 100a remains functional due to the failure of the other hosts
21 100a in the primary group 130a, and, as a result, the one
22 remaining functional host 100a is unable to perform assigned
23 tasks by itself. In this instance, the application 103b will run
24 for the system 50 to perform the assigned tasks.

1 In one embodiment, the heartbeat check 101a and heartbeat
2 check 101b are modules, software programs, firmware, hardware, a
3 combination of these components, or other suitable components.

4 The clustering programs 104a and 104b permit the hosts 100a
5 and 100b to function as a cluster computing system and are
6 conventionally known programs. The heartbeat check 101a can be
7 separate from the clustering program 104a, or may be combined or
8 attached with the clustering program 104a as one program.

9 Similarly, the heartbeat check 101b can be separate from the
10 clustering program 104b, or may be combined or attached with the
11 clustering program 104b as one program.

12 The operating system 102a provides APIs (application
13 program interfaces (APIs) for the Clustering Program 104a and
14 the Heartbeat Check 101 to use. For example, the operating
15 system 102a provides "open", "read", "write", and "close" to the
16 storage volumes. Heartbeat Check 101 uses these APIs when,
17 e.g., sending a heartbeat message (e.g., "open(vol)" to get a
18 pointer to the volume, "write(message)" to write a message, and
19 "close (vol)" to discard the pointer).

20 The paths 120a and 120b in Figure 1 transmit information
21 between the hosts 100a and the storage system 110a by use of a
22 standard protocol. Examples of the path 120 include SCSI, Fibre
23 channel, ESCON, or Ethernet, which standard protocols are SCSI-
24 3, FCP, ESCON, and TCP-IP, respectively.

1 Similarly, an operating system 102b performs functions for
2 hosts 100b, as similarly described above with the functions of
3 operating system 102a.

4 Each host has a Clustering program 104, heartbeat check
5 101, and operating system 102a. The heartbeat check 101 may be
6 a part of clustering program 104 (not separated). Each
7 operating system 102a works independently. Clustering program
8 104 (and heartbeat check 101) know the state of other hosts 100
9 (i.e., if the host is dead or alive). Based on the detected
10 state of the host, the clustering program determines to perform
11 or not perform fail-over.

12 Each host 100a has its own application 103a if a user
13 specifies accordingly. For example, a Host1 (in hosts 100a) may
14 run an Oracle database, Host2 may run a payroll application,
15 Host3 may run an order entry application, and the like. If
16 Host1 fails, then Oracle database is opened at, e.g., Host2.
17 Thus, Host2 now runs the Oracle database and the payroll
18 application.

19 The present invention chooses one host in the primary group
20 130a as a master host 160a, and one host in the secondary group
21 130b as a master host 160b. As described below, the master hosts
22 160a and 160b send "heartbeat" signals 300 to each other to
23 determine if a fail over should be performed. Other hosts 100a
24 in primary group 130a may become a master host 160a if the

1 current master host 160a is deemed to fail in following some
2 rules. Some examples of rules may include the following:

3 (1) The master host 160a did not send a heartbeat message
4 for 1 minute; or

(2) The master host 160a sent an invalid message (e.g., the message include an invalid date and time (not current), an invalid ID of the master host, and/or an expired instance (or process) ID of the cluster).

9 Similarly, other hosts 100b may become a master host 160b in
10 secondary group 130b if the current master host 160b is deemed to
11 fail in following some rules.

12 All of the hosts 100a (including master host 160a) are
13 connected by a network 140 to all of the hosts 100b (including
14 master host 160b). Thus, any of the hosts 100a in primary group
15 130a can communicate with any of the hosts 100b in the secondary
16 group 130b. Typically, the network 140 may be a Local Area
17 Network (LAN) or a Wide Area Network (WAN) such as the Internet.

As known to those skilled in the art, a LAN is a computer network that spans a relatively small area. Most LANs are confined to a single building or group of buildings. Most LANs connect workstations and personal computers (PCs). Each node (individual computer) in a LAN has its own central processing unit (CPU) with which it executes programs, but it is also able to access data and devices anywhere on the LAN. Thus, many

1 users can share expensive devices, such as laser printers, as
2 well as data. Users can also use the LAN to communicate with
3 each other, by, for example, sending e-mail or engaging in chat
4 sessions. There are many different types of LANs, with
5 Ethernets being the most common for PCs. LANs are capable of
6 transmitting data at very fast rates, much faster than the data
7 transmitted over a telephone line. However, the distances over
8 LANs are limited, and there is also a limit on the number of
9 computers that can be attached to a single LAN.

0 As also known to those skilled in the art, a WAN is a
1 computer network that spans a relatively large geographical
2 area. Typically, a WAN includes two or more LANs. Computers
3 connected to a WAN are often connected through public networks,
4 such as the telephone system. They can also be connected
5 through leased lines or satellites. The largest WAN in
6 existence is the Internet.

7 Through the network 140, the master hosts 160a and 160b can
8 send a heartbeat signal to each other. The network 140 also
9 permits master host 160a and 160b to perform heartbeat checking
10 with each other. In other words, the master host 160a can check
11 whether the master host 160b is alive (functional) or not by
12 checking for a heartbeat signal from the master host 160b, as
13 described below. Similarly, the master host 160b can check

1 whether the master host 160a is alive or not by checking for a
2 heartbeat signal from the master host 160a, as described below.

3 The primary group 130a is coupled to a storage system 110a
4 in the production site, and the secondary group 130b is coupled
5 to a storage system 110b in the standby site. Each of the
6 storage systems 110a and 110b form, for example, a disk system.
7 Each of the storage systems 110a and 110b comprises two or more
8 disks. The storage systems 110a and 110b are connected to each
9 other by one or more remote links 150 through which the storage
10 systems 110a and 110b communicate with each other. Typically,
11 the remote links 150 may be ESCON, Fibre Channel,
12 telecommunications lines, or a combination that may include
13 ESCON, Fibre Channel, and telecommunication lines.

14 As known to those skilled in the art, ESCON is a set of
15 products, e.g., from IBM, that interconnect S/390 computers with
16 each other and with attached storage, locally attached
17 workstations, and other devices using optical fiber technology
18 and dynamically modifiable switches called ESCON Directors.

19 As also known to those skilled in the art, Fibre Channel is
20 a serial data transfer architecture developed by a consortium of
21 computer and mass storage device manufacturers and now being
22 standardized by the American National Standards Institute
23 (ANSI). The most prominent Fibre Channel standard is the Fibre
24 Channel Arbitrated Loop (FC-AL) which is designed for new mass

1 storage devices and other peripheral devices that require very
2 high bandwidth. Using optical fiber to connect devices, FC-AL
3 supports full-duplex data transfer rates of 100 megabytes per
4 second (MBps).

5 The disk system (formed by storage systems 110a and 110b)
6 forms a remote data mirroring system and comprises one or more
7 remote mirror 111. Each remote mirror 111 comprises a storage
8 volume (heartbeat PVOL) 111a in storage system 110a and a
9 storage volume (heartbeat SVOL) 111b in storage system 110b.

10 When the heartbeat check 101a writes a heartbeat message 300 to
11 the heartbeat PVOL 111a, the storage system 110a then writes the
12 heartbeat message 300 to the heartbeat SVOL 111b via remote link
13 150. The heartbeat check 101b then reads the heartbeat signal
14 300 from the heartbeat SVOL 111b to check if the hosts 100a are
15 alive.

16 The number of remote mirrors 111, heartbeat PVOLs 111a,
17 heartbeat SVOLs 111b, and remote links 150 (linking a heartbeat
18 PVOL 111a with a heartbeat SVOL 111b) may vary. A heartbeat
19 PVOLs 111a may fail due to hardware problems. The use of two or
20 more heartbeat PVOLs 111a for use by heartbeat signals 300
21 advantageously achieves higher reliability for the system 50.

22 The heartbeat check 101a writes the heartbeat message 300
23 via path 170a to the heartbeat PVOL(s) 111a by use of, for
24 example, a Small Computer System Interface (SCSI) driver. SCSI

1 is a parallel interface standard used by Apple Macintosh
2 computers, PCs, and many UNIX systems for attaching peripheral
3 devices to computers. SCSI interfaces provide for faster data
4 transmission rates (up to about 80 megabytes per second) than
5 standard serial and parallel ports. The storage system 110a
6 sees the heartbeat signal 300 as a write data. The storage
7 system 100a stores the heartbeat signal 300 in heartbeat PVOL(s)
8 111a and also transmit the heartbeat signal 300 along remote
9 link 150 by use of a conventional driver or transceiver (not
0 shown in Figure 1).

1 The heartbeat signal 300 is received by a conventional
2 remote copy mechanisms in the storage system 110b and written to
3 the SVOL(s) 111b. The heartbeat check 101b then reads the
4 heartbeat signal 300 data stored in the heartbeat SVOL(s) 111b
5 via path 170b by use of conventional APIs that the operating
6 system provides.

7 The heartbeat check 101a (in master server 160a) sends a
8 heartbeat signal 300 at pre-determined intervals such as, for
9 example, every second, every 10 seconds, every 60 seconds, etc.
0 The heartbeat check 101a sends a heartbeat signal 300 along the
1 path 170a to one or more heartbeat PVOL(s) 111a as described
2 above. In addition to being able to send heartbeat signals 300
3 via remote link 150, the heartbeat check 101a can also

FIG. 10

1 simultaneously send heartbeat signals 300 along network 140 to
2 the hosts 100b.

3 The disk system (formed by storage systems 110a and 110b)
4 further comprises one or more remote mirror 112 for storing
5 production data. Each remote mirror 112 comprises a storage
6 volume (user's PVOL 112a) and a storage volume (user's SVOL
7 112b). As an example, a user's PVOL 112a or 112b comprises a
8 database such as a database available from Oracle Corporation.

9 The user's PVOL 112a or 112b may also be storage volumes for
10 storing data from the World Wide Web, text files, and the like.
11 When the application 103a updates data on the user's PVOL 112a,
12 the storage system 110a writes the data to the user's SVOL 112b
13 by use of a conventional remote copy mechanism to transmit the
14 data across the remote link 151 to storage system 110b. The
15 storage system 110b receives the data transmitted along remote
16 link 151 by use of a conventional remote copy mechanisms, and
17 the received data is then written into the user's SVOL 112b.

18 Hosts 100b (including master host 160b) access the user's
19 SVOL(s) 112b to read stored data after fail-over to secondary
20 group 130b occurs. In other words, if failure occurs in the
21 production site so that the primary group 130a is unable to
22 perform assigned operations or tasks, then the hosts 100b in the
23 secondary group 130b in the standby site will perform the
24 operations and tasks for system 50. Examples of failures that

1 may trigger a fail-over includes host failure, storage system or
2 disk failure, applications or software failure, hardware
3 failure, signal paths or connections failure, and other types of
4 failures in the production site that will prevent the host group
5 130a from performing assigned operations or tasks for system 50.

6 As known to those skilled, mirrored volumes containing
7 production data (e.g. user database) are sometimes broken
8 manually by having the user issue a break (split) command. The
9 mirrored volumes are broken for purposes of, for example,
10 performing backup tasks, running other applications or jobs on
11 the user's SVOL(s) 112b, and the like. A user (in the
12 production site) may issue a split (break) command, to prevent
13 the user's PVOL 112a from sending via remote link 151 data to
14 user's SVOL 112b. In other words, a split command prevents the
15 storage system 110a from sending data to the storage system
16 110b.

17 Thus, the remote mirror 111 is separate from the remote
18 mirror 112 because the remote mirror 111 should not be subject
19 to a split command. As stated above, the heartbeat PVOL 111a
20 sends via remote link 150 the heartbeat signals 300 to the
21 heartbeat SVOL 111b. A split command would disadvantageously
22 prevent the heartbeat PVOL 111a from sending via remote link 150
23 the heartbeat signals 300 to the heartbeat SVOL 111b. By
24 separating the remote mirrors 111 and 112, the heartbeat PVOL

1 111a can continue to send the heartbeat signals 300 to the
2 heartbeat SVOL 111b via remote link 150 even if the user issues
3 a split command that prevents the user's PVOL 112a from sending
4 data to the user's SVOL 112b via remote link 151.

5 This invention prevents the user from using mirrors 112
6 containing production data for the heartbeat checking tasks, and
7 may also alert the user if he/she tries to break the mirrored
8 volumes before de-installing the mirrored volumes 111 for
9 heartbeat signals.

10
11 TABLES

12 Heartbeat Status Table 250: Figure 2 is a block diagram of
13 a Heartbeat Status Table 250 stored in each of the master hosts
14 160a and 160b. The Heartbeat Status Table 250 is stored in both
15 the memory in a host and a volume. The volume prefers to be
16 remotely mirrored (so is the PVOL). The table 250 is used by
17 Heartbeat check 101a. The heartbeat check 101a running on
18 master host 160a will create, refer to, and change a Heartbeat
19 Status Table 250. Similarly, the heartbeat check 101b running
20 on master host 160b will create, refer to, and change another
21 Heartbeat Status Table 250 in the master host 160b. When
22 another one of the hosts 100a becomes the master host 160a, the
23 hosts 100a will also create, refer to, and change an associated
24 Heartbeat Status Table 250. Similarly, when another one of the

1 hosts 100b becomes the master host 160b, the hosts 100b will
2 also create, refer to, and change an associated Heartbeat Status
3 Table 250.

4 As discussed in detail below, the Heartbeat Status Table
5 250 comprises: Network Heartbeat Enable 200, Remote Group Status
6 210, Remote Copy Heartbeat Enable 220, Remote Group Status 230,
7 Device Addresses (1), (2),... (n) 240, and Device Status (1),
8 (2),... (n) 241.

9 The Network Heartbeat Enable 200 shows whether the network
10 140 can be used for transmitting a heartbeat signal 300.

11 Possible values include the following: "ENABLE", "DISABLE", and
12 "FAILED". The entry "ENABLE" indicates that the user has
13 enabled the system 50 to permit a heartbeat signal 300 to
14 transmit across network 140. When user specifies not to use the
15 network 140 for heartbeat, Network Heartbeat Enable 200 turns
16 this entry into "DISABLE". The user can permit or disable the
17 sending of heartbeat signals 300 across the network 140 by
18 issuing commands in the master host 160a. When the heartbeat
19 check 101a (or 101b) finds unrecoverable errors in system 50,
20 then the heartbeat check 101a turns the Network Heartbeat Enable
21 200 entry into "FAILED". The heartbeat check 101a (or 101b)
22 does not use the network 140 for checking heartbeat signals if
23 the Network Heartbeat Enable 200 shows the entries of "DISABLE"
24 or "FAILED".

FOR THE "5096021"

1 will show "FAILED". As discussed below, the entry in the Device
2 Status 241 will show whether each of the remote mirrors 111 are
3 enabled or disabled.

4 The results of performing a heartbeat check via a remote
5 mirror 111 shows whether the remote group 130b is alive
6 (functional) or non-functional. The Remote Group Status 230
7 shows the results of the heartbeat check. The status of this
8 entry in the Remote Group Status 230 depends only on performing
9 the heartbeat checking via remote mirrors 111. If all of the
10 remote links 150 fail or all remote mirrors 111 fail, then a
11 heartbeat signal 300 via remote mirrors 111 cannot reach the
12 remote group 130b from the primary group 130a. As a result, the
13 entry in the Remote Group Status 230 will show "FAILED".

14 Device Address 240 shows a device address of a mirror 111
15 for a heartbeat signal 300. For example, an entry in device
16 address 240 is stored in master host 160a and contains a device
17 address of a heartbeat PVOL 111a. The heartbeat check 101a will
18 write heartbeat signals 300 to the heartbeat PVOLs 111a in a
19 remote mirror 111 with an addresses listed in device address 240
20 of the associated heartbeat status table 250 in the master host
21 160a. This same type of entry in the device address 240 is also
22 stored in master host 160b and contains the same device address
23 of the remote mirror 111 with the heartbeat SVOL 111b. The
24 heartbeat check 101b will read heartbeat signals 300 stored in

1 the heartbeat SVOLs 111b in the mirror 111 with the address.
2 listed in the device address 240.

3 Device Status 241 shows a status of the device (heartbeat
4 mirror 111) that is registered with Device Address 240. The
5 value in the Device Status 241 entry may include "ENABLE",
6 "DISABLE", or "FAILED". When a user deactivates a particular
7 mirror 111 or a failure occurs in that particular mirror 111,
8 the Device Status 241 of that particular mirror 111 will show
9 "DISABLE" or "FAILED", and the heartbeat check 101a does not use
10 that particular failed mirror for the heartbeat signal 300
11 transmission. As shown in Figure 2, there are two or more
12 entries (240a, 240b,... 240c) in the Device Address 240 and two or
13 more entries (241a, 241b,... 241c) in the Device Status 241 for
14 two or more mirrors 111 for processing the heartbeat signal 300.
15 In other words, entries 240a and 241a are associated with a
16 mirror 111, while entries 240b and 241b are associated with
17 another mirror 111. Unused entries in the Heartbeat Status
18 Table 250 contain "NULL".

19

20 Heartbeat message 300

21 The master host 160a at the primary site sends a heartbeat
22 message (or signal) 300 to the master host 160b at the standby
23 site through the network 140 and/or through the heartbeat mirror
24 111. When using the mirror 111 to send the heartbeat messages

1 300, the master host 160a writes this heartbeat message 300 to
2 the heartbeat PVOL 111a, and the master host 160b reads the
3 transmitted heartbeat message 300 from the heartbeat SVOL 111b.

4 Figure 3 is a block diagram showing an example of the data
5 format of a heartbeat message. A heartbeat message 300 includes
6 at least some of the following entries. A Serial Number 310 is
7 a number serially assigned to the heartbeat message 300. In one
8 embodiment, this number increments (or counts up) by one (1) for
9 each heartbeat message 300 that is sent, and the value of the
10 serial number is re-initialized to 1 for the next heartbeat
11 message 300 that is sent after a heartbeat message with the
12 maximum value of the serial number 310 is sent. A Time 320
13 contains the time when the Heartbeat check 101a running on a
14 master host 160a generates a heartbeat message 300. An
15 identifier 330 is used to identify the sender of the message.
16 This identity may be, for example, a unique number assigned to
17 the heartbeat check 101a running at the primary site, a name
18 uniquely assigned to the heartbeat check 101a that is sending
19 the heartbeat message 300, the Internet Protocol (IP) address of
20 the master host 160a, or combination of the above
21 identifications.

22
23 Method of installing mirrors used for heartbeat signals (see
24 Figure 4)
25

SECRET

1 As stated above, the heartbeat mirrors 111 used for
2 heartbeat signals 300 are different from the production data
3 mirrors 112 used for storing the production data. This is
4 because the production mirrors 112 may be broken manually for
5 other purposes, such as performing a backup from SVOL, running
6 other jobs or applications on SVOL, etc.

7 Figure 4 is a flowchart diagram illustrating a method of
8 installing a heartbeat mirror 111 used for heartbeat signals
9 300. The heartbeat check 101a provides a user interface for
10 creating 400 heartbeat mirrors 111 used for the transmission and
11 storage of heartbeat signals 300. The heartbeat check 101a also
12 provides the user interface to control the mirrors 111, such as,
13 for example, the creation, deletion, and the breaking
14 (splitting) of mirrors 111. As previously noted, the mirrors
15 112 (which contain production data) are not used for heartbeat
16 signal 300 transmission and processing. The heartbeat check
17 101a provides the user interface to activate and/or deactivate
18 mirrors 111 for heartbeat signal 300 transmission or processing.
19 Using this user interface, a user can activate 410 any or all of
20 the heartbeat mirrors 111. The input parameters for the user
21 interface are, for example, the heartbeat PVOL device 111a
22 address and the heartbeat SVOL device 111b address. In this
23 step 410, if a mirror 111 has not been created, then the

1 heartbeat check 101a halts the process and displays on the user
2 interface an alert message, such as "no mirror is created".

3 Once the user activates 410 the heartbeat mirror 111, in
4 step 420 the heartbeat check 101a running on the master host
5 160a registers the heartbeat PVOL 111a (in activated mirror 111)
6 to the Device Address 240 (Figure 2), and changes Device Status
7 241 to "ENABLE" in Heartbeat Status Table 250 (Figure 2).

8 Before performing this step 420, the heartbeat check 101a
9 displays an alert message in the user interface that production
10 data should not be placed in the mirrors 111 used for the
11 heartbeat signals 300. The heartbeat check 101a then sends 430
12 to the heartbeat check 101b running on the standby site the
13 following: an activation message along with the parameters.
14 These parameters include the address of the SVOL 111b in the
15 mirror 111 to be activated. The heartbeat check 101a sends this
16 activation message via network 140 or by using a heartbeat
17 mirror 111 that is already available. When the heartbeat check
18 101b, running on the standby site, receives the activation
19 message sent by heartbeat check 101a in step 430, then in step
20 440 the heartbeat check 101b registers the heartbeat SVOL(s)
21 111b to the Device Address 240 and changes the Device Status 241
22 to "ENABLE" in Heartbeat Status Table 250. Thus, the heartbeat
23 mirror 111 is now installed, and the heartbeat check 101a can
24 now send heartbeat signals 101b to the heartbeat check 101b via

1 remote link 150. It is noted that a plurality of mirrors may be
2 installed when the method shown in Figure 4 is performed.

3

4 Method of de-installing mirrors 111 used for heartbeat messages
5 300 (see Figure 5)
6

7 The user may want to de-install a mirror 111 that is being
8 used for the heartbeat messages 300. For example, the user may
9 decide to decrease the number of mirrors 111 used for heartbeat
10 messages 300 since the performance of heartbeat checking
11 degrades if many mirrors 111 are used for transmitting and
12 processing heartbeat messages 300.

13 Figure 5 is a flowchart diagram of a method of de-installing
14 a mirror 111 used for heartbeat messages 300. The user de-
15 activates a mirror (or mirrors) 111 by using the user interface
16 provided by the heartbeat check 101a. The heartbeat check 101a,
17 running at the production site, de-activates 500 the particular
18 mirror(s) 111 specified by the user. To de-activate the
19 particular mirror(s) 111 specified by the user, the heartbeat
20 check 101a changes 510 the entries in the Device Address 240 and
21 Device Status 241 in the Heartbeat Status Table 250 (Figure 2)
22 to "NULL". The heartbeat check 101a then sends 520 a de-
23 activation message along with the parameters to the heartbeat
24 check 101b running at the standby site. These parameters are
25 device address of SVOL(s). The de-activation message is sent is

09760345-01201
T02270

1 via network 140 or by using any mirror 111 that is still
2 available for transmitting signals. The heartbeat check 101b,
3 running at the standby site, de-activates the particular
4 mirror(s) 111 specified in step 500 by the user. To de-activate
5 the particular mirror(s) 111, the heartbeat check 101b changes
6 530 the associated entries in the Device Address 240 and Device
7 Status 241 in the Heartbeat Status Table 250 to "NULL". If the
8 user no longer needs to use the de-activated mirror(s) 111,
9 he/she can delete 540 the mirror(s) 111 using a known user
10 interface provided by storage system vendors. Deactivating a
11 mirror prohibits the Heartbeat Check to use the mirror for
12 sending a heartbeat message, but the mirror is still formed (the
13 PVOL and SVOL relationship is still valid). Deleting a mirror
14 actually deletes the mirror. Thus, the mirror relationship no
15 longer exists. A User may need to delete mirrors for some
16 reasons (e.g., to increase performance and the like). In such
17 cases, the user should deactivate the mirrors before deleting
18 the mirrors, in order for the Heartbeat Check to not try to send
19 a message via the deleted mirrors.

20 It is also noted that when the user tries to break (split)
21 a heartbeat mirror(s) 111 and if the particular heartbeat
22 mirror(s) 111 is not de-installed, then the heartbeat check 101a
23 will not break the particular heartbeat mirror(s) 111 until the

1 mirror(s) is de-installed. This insures that an installed
2 heartbeat mirror 111 will not be split.

3

4 Method of Sending a Heartbeat message 300 (see Figure 6)

5 Figure 6 is a flowchart diagram of a method of sending a
6 heartbeat message 300 in accordance with an embodiment of the
7 present invention. The heartbeat check 101a sends periodically,
8 for example, every one-minute, a heartbeat message 300 to the
9 heartbeat check 101b as shown in this flowchart. The user can
10 specify the interval during which a heartbeat message 300 is
11 transmitted.

12 The heartbeat check 101a first determines 600 whether the
13 network 140 can be used for transmitting a heartbeat signal 300.
14 The Network Heartbeat Enable 200 (Figure 2) entry in the
15 Heartbeat Status Table 250 shows whether the network 140 can be
16 used for transmitting a heartbeat signal 300. If the network
17 140 can be used for transmitting a heartbeat signal 300, then
18 the Network Heartbeat Enable 200 entry will indicate "ENABLE".
19 In this case, the heartbeat check 101a sends 610 a heartbeat
20 message 300 via network 140. To create a heartbeat message 300,
21 the heartbeat check 101a increments the current value of a
22 Serial Number 310 (Figure 3), obtains the current time from a
23 timer in the operating system 102a, and places these information

09760345-01201
T0210-01201

1 into the heartbeat message 300 along with a predetermined
2 identifier 330.

3 If, in step 600, the network 140 can not be used for
4 transmitting heartbeat signals 300 (i.e., the Network Heartbeat
5 Enable 200 entry does not indicated "ENABLE"), then the method
6 proceeds to step 620 which is discussed below.

7 As stated above, the Remote Copy Heartbeat Enable 220 entry
8 in the Heartbeat Status Table 250 (Figure 2) shows if the remote
9 mirrors 111 used for heartbeat signals 300 are available. When
10 the user specifies to use all the remote mirrors 111 for
11 heartbeat signals 300, then the entry in the Remote Copy
12 Heartbeat Enable 220 will show "ENABLE". When the user
13 specifies not to use all of the remote mirrors 111 for heartbeat
14 signals 300, then the entry in the Remote Copy Heartbeat Enable
15 220 will show "DISABLE". When all remote mirrors 111 are
16 disabled or have failed, then the entry in the Remote Copy
17 Heartbeat Enable 220 will show "FAILED".

18 The heartbeat check 101a checks 620 if the Remote Copy
19 Heartbeat Enable 220 entry shows "ENABLE" (i.e., at least one
20 remote mirror 111 is available for use by the heartbeat messages
21 300. If so, then the heartbeat check 101a sends the heartbeat
22 message 300 via remote link 150. If the Remote Copy Heartbeat
23 Enable 220 entry does not show "ENABLE" (i.e., all remote
24 mirrors are unavailable for use by the heartbeat messages 300),

1 then the method ends. Thus, "ENABLE" shows at least one remote
 2 mirror is available, and "DISABLE" shows all remote mirrors are
 3 unavailable. A message cannot be sent when "DISABLE" is shown.

4 The heartbeat check 101a writes 630 a heartbeat message 300
 5 to all available heartbeat PVOLs 111a. To check if a heartbeat
 6 PVOL(s) 111a is available, the heartbeat check 101a checks the
 7 status of every entry of Device Status 241 in Heartbeat Status
 8 Table 250 to determine all mirrors 111 that are available. The
 9 heartbeat check 101a then writes the heartbeat message 300 to
 10 the heartbeat PVOL devices 111a in the available mirrors 111.

11 As stated above, a mirror 111 is available if its associated
 12 entry show "ENABLE" in the Device Status 241. The heartbeat
 13 check 101a does not send heartbeat signals 300 to the heartbeat
 14 PVOL devices 111a that are in mirrors 111 that are unavailable.
 15 As stated above, a heartbeat mirror 111 is unavailable if it has
 16 an associated entry of "NULL" in the Device Status 241.

17 If a network 140 failure occurs while the heartbeat check
 18 101a is sending a heartbeat message 300, then the Network
 19 Heartbeat Enable 200 entry will indicate "FAILED". If a device
 20 (heartbeat mirror 111) failure occurs while the heartbeat check
 21 101a is writing a heartbeat message 300 to a heartbeat mirror
 22 111, then the Device Status 241 entry for that failed mirror
 23 will indicate "FAILED". At this time, the heartbeat check 101a
 24 checks the other entries of Device Status 241 (i.e., the

1 heartbeat check 101a checks if other heartbeat mirrors 111 are
2 available so it can determine which other heartbeat PVOLs 111a
3 may be used for the processing of the heartbeat signals 300).
4 The heartbeat check 101a will indicate the entry in the Remote
5 Copy Heartbeat Enable 220 as "FAILED" if all the entries of
6 Device Status 241 show either "FAILED", "DISABLE", or "NULL".
7

8 Method of receiving a heartbeat message 300 (see Figure 7)

9 The heartbeat check 101b (of master host 160b) periodically
10 receives and checks for heartbeat messages 300 sent by the
11 heartbeat check 101a. If there is one or more heartbeat mirrors
12 111 that is functioning, then the heartbeat check 101b reads
13 from each heartbeat SVOL 111b (in each functioning heartbeat
14 mirror 111) until the heartbeat check 101b finds a valid
15 heartbeat message 300 stored in at least one of the heartbeat
16 SVOLs 111b.

17 The definition of the valid heartbeat message 300 includes
18 one or more of the following:

19 (1) Based on the Identifier 330 in a heartbeat message 300,
20 the heartbeat check 101b approves the heartbeat message 300 sent
21 from the Heartbeat check 101a;

22 (2) The Serial Number 310 is continuously incremented
23 within a timeout period (e.g., one minute); and

1 (3) The Time 320 is continuously updated within a timeout
2 period.

3 Other definitions of a valid heartbeat may also be made, as
4 specified by the user.

5 It is noted further that the above condition (1) specifies
6 that the sender of a message is a member (host) of a cluster.
7 Within a cluster, each host knows and can identify the members
8 of the cluster. Thus, the receiver of a message can identify
9 whether the message is sent from a member of the cluster.

10 It is noted further that in the above condition (2), a
11 receiver observes messages sent by a sender. If the serial
12 number 310 is not incremented within, say, one minute, then the
13 sender is deemed as failed.

14 It is noted further that in the above condition (3), a
15 receiver observes messages sent by a sender. If Time 320 is not
16 updated within, say, one minute, then the sender is deemed as
17 failed.

18 Referring now to Figure 7, there is shown a flowchart
19 diagram of a method of receiving a heartbeat message 300 in
20 accordance with an embodiment of the present invention. The
21 heartbeat check 101b checks 700 if the network 140 can be used
22 for transmission of heartbeat signals 300 (i.e., if the Network
23 Heartbeat Enable 200 entry in Heartbeat Status Table 250 shows
24 "ENABLE"). If the network 140 can be used for heartbeat signals

1 300, then the heartbeat check 101b checks 710 if a valid
 2 heartbeat message 300 has been received via network 140. If a
 3 valid heartbeat message 300 has not been received, then the
 4 heartbeat check 101b skips its checking of heartbeat signals 300
 5 received via network 140 and marks the network 140 and
 6 production group 130a as having failed by changing the entries
 7 in Network Heartbeat Enable 200 and Remote Group Status 210 as
 8 "FAILED".. This indicates that after the heartbeat check 101b
 9 checked the network 140 for the heartbeat signals 300, the
 10 production group 130a and the network 140 are regarded as having
 11 failed. As a result, the failed network 140 is not used for
 12 heartbeat checking operations by heartbeat check 101b.

13 The heartbeat check 101b checks 720 if the remote mirrors
 14 111 used for heartbeat signals 300 are available. If the Remote
 15 Mirror Heartbeat Enable 220 entry shows "ENABLE", then the
 16 heartbeat check 101b checks 730 for a received heartbeat message
 17 300 by checking at least one remote mirror 111 which is
 18 available. If none of the remote mirrors 111 are available,
 19 then the heartbeat check 101b skips the heartbeat checking
 20 operation via remote mirror(s) 111 and proceeds to step 740
 21 which is discussed below.

22 The heartbeat check 101b reads 730 a heartbeat message 300
 23 from each heartbeat SVOL 111b. If the heartbeat check 101b
 24 finds a valid heartbeat message 300, then the standby host group

1 130b will not take over operations for the production group
2 130a, since the production group 130a is deemed as alive (has
3 not failed). On the other hand, if the heartbeat check 101b
4 finds all the heartbeat messages 300 as invalid, then the
5 heartbeat check 101b marks the entries in the Remote Copy
6 Heartbeat Enable 220 and Remote Group Status 230 as "FAILED".
7 This indicates that production group 130a and all remote mirrors
8 111 are regarded as having failed based upon the result of the
9 heartbeat checking via remote mirroring, and all remote mirrors
10 111 are not used for heartbeat checking from then on.

11 If the heartbeat check 101b finds a particular remote
12 mirror 111 that contains an invalid heartbeat message 300, then
13 the heartbeat check 101b marks the Device Status 241 entry
14 associated with that particular remote mirror 111 as "FAILED".
15 As a result, the heartbeat SVOL 111b in that remote mirror 111
16 is not used for process of heartbeat checking by use of the
17 remote mirrors.

18 After the above steps has been performed, the Heartbeat
19 Status Table 250 (Figure 2) will contain the results of the
20 heartbeat checking via network 140 and heartbeat checking by use
21 of remote mirroring. If neither the Remote Group Status 210 nor
22 the Remote Group Status 230 shows the entry "ALIVE", then
23 Heartbeat check 101b regards production group 130a as dead and
24 will perform 740 the fail-over operation as described above. As

1 a result of the fail-over operation, the standby group 130b will
2 assume operation of the system 50 of Figure 1.

3 4 Non-stop addition and deletion of heartbeat mirrors 111

5 It is advantageous if the addition and deletion of a
6 heartbeat mirror(s) 111 are performed without affecting the
7 heartbeat checking operation described above. To achieve this
8 feature, the clustering system 50 (Figure 1) starts to use each
9 newly-created mirror volumes 111a and 111b in a new heartbeat
10 mirror 111. The system 50 stops in using existing mirror
11 volumes 111a and 111b in each heartbeat mirror 111 deleted by
12 the user.

13 As described in Figure 4, when installing a heartbeat
14 mirror 111 for use by heartbeat messages 300, the heartbeat
15 checks 101 (101a and 101b) register information of the newly-
16 created heartbeat mirror 111 to vacant entries (containing
17 "NULL") in Device Address 240 (e.g., Device Address 204c) and in
18 Device Status 241. On the other hand, as described above with
19 regard to Figure 6 and Figure 7, the heartbeat checks 101 (101a
20 and 101b) do not use the vacant entries in the Device Address
21 240 and Device Status 241. This is the same procedure when
22 deleting a heartbeat mirror 111. The heartbeat checks 101 will
23 stop using a heartbeat mirror 111 that is deleted or de-
24 installed. Thus, the heartbeat checks 101 do not need to be

1 stopped in its processing of heartbeat signals 300, while the
2 heartbeat checks 101 are adding (or are deleting) a heartbeat
3 mirror 111.

4
5 Method for setting the heartbeat checking procedure (Figure 8)

6 There are three methods for sending heartbeat messages 300
7 from the primary group 130a to the standby group 130b (or from
8 the standby group 130b to the primary group 130a). The
9 heartbeat messages 300 can be selectively sent: (1) through the
10 network 140, (2) through at least one remote mirrors 111, or (3)
11 through both the network 140 and at least one remote mirror 111.
12 The user can choose one of these methods for sending the
13 heartbeat messages 300. When the user indicates a method for
14 sending the heartbeat messages 300, the heartbeat check 101a
15 updates the Heartbeat Status Table 250 without affecting
16 heartbeat checking operation.

17 Changing the procedure for sending the heartbeat messages
18 300 is very useful when the network 140 or the remote mirrors
19 111 are being diagnosed, and or when regular maintenance work is
20 being performed.

21 Figure 8 is a flowchart diagram illustrating a method of
22 setting the heartbeat checking procedure in accordance with an
23 embodiment of the present invention. The user first requests
24 800 to change of the heartbeat checking procedure by indicating

1 if the heartbeat messages will be sent by use the network 140,
2 by use of the remote mirrors 111, or by use of both the network
3 140 and the remote mirrors 111. This request indicates that the
4 network 140 and the remote mirrors 111 for checking heartbeat
5 messages 300 are enabled or disabled.

6 Once a heartbeat checking procedure is indicated by the
7 user's request, the heartbeat check 101a and 101b will execute
8 the following. The user's request will activate (or deactivate)
9 810 the network 140 for heartbeat checking operations. The
10 value of the Network Heartbeat Enable 200 entry will be as
11 follows: If the user is activating the heartbeat checking via
12 network 140, then the entry in the Network Heartbeat Enable 200
13 will be "ENABLE". The heartbeat check 101a will send heartbeat
14 signals 300 along the network 140. If the user is deactivating
15 the heartbeat checking via network 140, then the entry in the
16 Network Heartbeat Enable 200 will be "DISABLE". The heartbeat
17 check 101a will not send heartbeat signals 300 along the network
18 140.

19 The user's request will activate (or deactivate) 820 a
20 remote mirror 111 for heartbeat checking operations. The value
21 of Remote Copy Heartbeat Enable 220 entry will be as follows:
22 If the user is activating the heartbeat checking via remote
23 mirror 111, then the entry in the Remote Copy Heartbeat Enable
24 220 will be "ENABLE". The heartbeat check 101a will send

1 heartbeat signals 300 via remote mirrors 111. If the user is
2 deactivating the heartbeat checking via remote mirror 111, then
3 the entry in the Remote Copy Heartbeat Enable 220 will be
4 "DISABLE". The heartbeat check 101a will not send heartbeat
5 signals 300 via remote mirrors 111.

6 The heartbeat check 101a sends 830 the user's request (made
7 in step 800) to the heartbeat check 101b. This request is sent
8 via network 140 or a remote mirror 111 that is currently
9 available.

10 The heartbeat check 101b then performs 840 and 850 a
11 similar process described in steps 810 and 820. Specifically,
12 The user's request will activate (or deactivate) 840 the network
13 140 for heartbeat checking operations. The value of the Network
14 Heartbeat Enable 200 entry will be as follows: If the user is
15 activating the heartbeat checking via network 140, then the
16 entry in the Network Heartbeat Enable 200 will be "ENABLE". The
17 heartbeat check 101b can check for heartbeat signals 300 along
18 the network 140. If the user is deactivating the heartbeat
19 checking via network 140, then the entry in the Network
20 Heartbeat Enable 200 will be "DISABLE". The heartbeat check
21 101b will not be able to check for heartbeat signals 300 along
22 the network 140.

23 The user's request will activate (or deactivate) 850 a
24 remote mirror 111 for heartbeat checking operations. The value

1 of Remote Copy Heartbeat Enable 220 entry will be as follows:
2 If the user is activating the heartbeat checking via remote
3 mirror 111, then the entry in the Remote Copy Heartbeat Enable
4 220 will be "ENABLE". The heartbeat check 101ab can check for
5 heartbeat signals 300 via remote mirrors 111. If the user is
6 deactivating the heartbeat checking via remote mirror 111, then
7 the entry in the Remote Copy Heartbeat Enable 220 will be
8 "DISABLE". The heartbeat check 101b will not be able to check
9 for heartbeat signals 300 via remote mirrors 111.

10 It is also possible to activate or deactivate a remote
11 mirror 111, or a set of remote mirrors 111 (i.e., not all remote
12 mirrors 111). To do this, the heartbeat check 101a (in step
13 820) and heartbeat check 101b (in step 850) changes the entry in
14 the Device Status 241 for an associated remote mirror(s) 111 to
15 "ENABLE" (if activating the mirror 111)" or to "DISABLE" (if
16 deactivating the remote mirror 111).

17 18 Bi-directional heartbeat messages (see Figure 9)

19 Reference is now made to Figure 9 which illustrates a block
20 diagram of a system 900 in accordance with another embodiment of
21 the present invention. A cluster system may require the
22 checking of heartbeat signals in a bi-directional manner. Thus,
23 the production group 130a may want to know whether the standby
24 group 130b is available or not available. For example, the user

1 at the production site may want to check the availability of the
2 standby group 130b. Thus, a bi-directional heartbeat mechanism
3 would be useful in this instance. For this bi-directional
4 mechanism in accordance with an embodiment of the present
5 invention, other mirrored volumes are created: a heartbeat PVOL
6 113a at the standby site and a heartbeat SVOL 113b at the
7 production site. The heartbeat PVOL 113a and heartbeat SVOL
8 113b are in the mirror 113. The master host 160b in the standby
9 group 130b writes a heartbeat signal 300' to the heartbeat
10 PVOL(s) 113a and the master host 160a in the production group
11 reads the heartbeat signal 300' from the heartbeat SVOL(s) 113b
12 to check if the standby group 130b is alive.

13 In this embodiment, the heartbeat check 101a not only sends
14 heartbeat messages 300 but also receives heartbeat messages 300'
15 from the heartbeat check 101b to check if the heartbeat check
16 101a is alive. To implement this embodiment, as shown in Figure
17 9, the remote mirror 113 is created, where the heartbeat PVOL
18 113a is in the storage system 110b at standby site, and the
19 heartbeat SVOL 113b is in the storage system 110a at production
20 site. The number of remote mirrors 113, heartbeat PVOLs 113a,
21 heartbeat SVOLs 113b, and remote links 150' (linking a heartbeat
22 PVOL 113a with a heartbeat SVOL 113b) may vary.

23 For the system 900 shown in Figure 9, the user installs a
24 remote mirror 113 for transmitting heartbeat messages 300' from

1 the storage system 110b to the storage system 110a via remote
2 link 150'. The heartbeat check 101b writes a heartbeat signal
3 300' to the heartbeat PVOL 113a, and the storage system 110b
4 writes the heartbeat signal 300' to the heartbeat SVOL 113b via
5 the remote link 150'. The heartbeat check 101a can read the
6 heartbeat signal 300' from the heartbeat SVOL 113b.

7 All tables mentioned above may be used in this embodiment.
8 Additionally the mechanisms or methods such as the installing
9 and de-installing mirrors 113', the sending of heartbeat
10 messages 300', the receiving heartbeat messages 300', and the
11 setting a heartbeat checking procedure for heartbeat signals
12 300' are performed similarly to the methods relating the
13 heartbeat signals 300. For example, to carry out the functional
14 operations for heartbeat signals 300', the roles or functions of
15 the heartbeat check 101a and heartbeat check 101b are reversed
16 in Figures 4, 5, 6, 7, and 8.

17 Figure 10 is a flowchart diagram illustrating a method of
18 failure notification in accordance with an embodiment of the
19 present invention. This method can be performed by, for
20 example, the system 50 in Figure 1. A remote mirror(s) 111
21 and/or the network 140 are activated 1000 so that the primary
22 group 130a can selectively send a failure indication message
23 1100 (Figure 11) along the activated remote mirror(s) 111 and/or
24 along the network 140. The mirror(s) 111 may be activated in

1 the same manner as previously described above. A check 1005 for
 2 failure can be made in the host group 130a. Components in the
 3 host group 130a have uniquely assigned serial numbers for
 4 purposes of identification. An IP address will work for this
 5 purpose. For example, each server 100a has a unique
 6 identification number. The storage system 110a (or 110b) has a
 7 unique identification number. The storage volumes PVOLs 111a
 8 and PVOLs 112a also have uniquely assigned identification
 9 numbers and uniquely assigned addresses. If a component fails,
 10 then the heartbeat check 101a is configured to determine the
 11 component identification number of that failed component.

12 In one embodiment, the heartbeat message 1100 (Figure 11)
 13 includes the following information:

14 (1) Failed parts information 1105 (shows failed parts,
 15 such as a failed host, network, disk drive, others): This
 16 information may be the ASCII character code, or the number
 17 uniquely assigned to the parts.

18 (2) Level of failure 1120 (one of the values of "SYSTEM
 19 DOWN", "SERIOUS", "MODERATE", or "TEMPORALLY"): Based on this
 20 information, the alert action will differ. For example, if the
 21 level of failure is "SYSTEM DOWN" or "SERIOUS", then the level
 22 of failure is notified to the system manager by phone at any
 23 time. If the level of failure is "TEMPORALLY", then the
 24 information is just logged or recorded.

1 (3) Parts information 1110 describes detailed information
2 about the failed parts): For example, if the failed part is a
3 host, then the parts information shows the IP Address of the
4 failed host. If the failed part is a drive, then the parts
5 information shows the drive serial number.

6 The heartbeat check 101a is also configured to determine an
7 address of the failed component. For example, the heartbeat
8 check 101a will determine the address of a failed storage volume
9 PVOL 112a. Additionally, the heartbeat check 101a may record
10 the time of failure of a component.

11 The heartbeat check 101a then sends 1010 a failure
12 indication message 1100 (Figure 11) via remote mirror 111 and/or
13 via network 140 in a manner similar to the transmission of
14 heartbeat signals 300 as described above. After the failure
15 indication message 1100 is received 1015 by the heartbeat check
16 101b, then heartbeat check 101b can read the failure indication
17 message and display information in the master host 160b
18 interface concerning the failure in the primary group 130a.
19 This display information may include, for example, the identity
20 of the failed component, the address of the failed component,
21 and the time during which failure of the component was detected.

22 Figure 11 is a block diagram illustrating an example of a
23 format of a failure indication message 1100 in accordance with
24 an embodiment of the present invention. As stated above, the

1 failure indication message 1100 may include the unique component
2 identification number 1105 the failed component, parts
3 information (e.g., an address) 1110 of the failed component, and
4 the time 1115 during which the failure of the component was
5 detected.

6 It is also within the scope of the present invention to
7 implement a program or code that can be stored in an
8 electronically-readable medium to permit a computer to perform
9 any of the methods described above.

10 Thus, while the present invention has been described herein
11 with reference to particular embodiments thereof, a latitude of
12 modification, various changes and substitutions are intended in
13 the foregoing disclosure, and it will be appreciated that in some
14 instances some features of the invention will be employed without
15 a corresponding use of other features without departing from the
16 scope of the invention as set forth.